

1 **Fully-automated, semantic segmentation of whole-body ¹⁸F-FDG PET/CT images**
2 **based on data-centric artificial intelligence**

3

4 Lalith Kumar Shiyam Sundar^{1*}, Josef Yu^{1,2*}, Otto Muzik³, Oana C. Kulterer², Barbara Fueger², Daria Kifjak^{2,4},
5 Thomas Nakuz², Hyung Min Shin⁵, Annika Katharina Sima², Daniela Kitzmantl², Ramsey D Badawi⁶, Lorenzo
6 Nardo⁶, Simon R Cherry⁶, Benjamin A Spencer⁶, Marcus Hacker² and Thomas Beyer¹

7

8 ¹Quantitative Imaging and Medical Physics (QIMP) Team, Center for Medical Physics and Biomedical
9 Engineering, Medical University of Vienna, Austria

10 ²Department of Biomedical Imaging and Image-Guided Therapy, Medical University of Vienna, Vienna,
11 Austria

12 ³Department of Pediatrics, Wayne State University School of Medicine, Children's Hospital of Michigan,
13 Detroit, MI, USA

14 ⁴Department of Radiology, University of Massachusetts Chan Medical School/UMass Memorial Health
15 Care, Worcester, MA, USA

16 ⁵Division of General Surgery, Department of Surgery, Medical University of Vienna, Vienna, Austria

17 ⁶Department of Biomedical Engineering and Radiology, University of California-Davis, Davis, CA, USA

18 **Corresponding Author**

19 Otto Muzik, PhD

20 Wayne State University School of Medicine, Detroit, Michigan, USA

21 Telephone: +1-(734)-223-3506

22 otto@pet.wayne.edu

23 First authors contact information:

24 Lalith Kumar Shiyam Sundar, PhD and Josef Yu MD (*Equal contribution)

25 QIMP Team, Center for Medical Physics and Biomedical Engineering, Medical University of Vienna,

26 Vienna, Austria

27 Telephone: +43 1 40400 55450

28 Lalith.shiyamsundar@meduniwien.ac.at

29

30 **Running title:** PET/CT Multi-organ segmentation

31 Word count: 4624

32 Figures: 6

33

34 **Keywords:** multi-organ segmentation, total-body PET, systems medicine, artificial neural networks,

35 automated segmentation

36

37 **Funding**

38 This work was supported in parts by NIH Research grant R01CA29422.

39 **ABSTRACT**

40 We introduce Multi-Organ Objective Segmentation (MOOSE) software that generates subject-specific,
41 multi-organ segmentation using data-centric AI principles to facilitate high-throughput systemic
42 investigations of the human body via whole-body PET imaging. **Methods:** Image data from two PET/CT
43 systems (uEXPLORER and Siemens TruePoint TrueView) was used in training MOOSE. For non-cerebral
44 structures, 50 WB-CT images were used, 30 of which were acquired from healthy controls (HC, 14M/16F)
45 and 20 datasets were acquired from oncology patients (14M/6F). Non-cerebral tissues consisted of 13
46 abdominal organs, 20 bone segments, subcutaneous fat, visceral fat, psoas, and skeletal muscle. An expert
47 panel performed manual segmentation of all non-cerebral structures except for subcutaneous, visceral
48 fat, and skeletal muscle, which were semi-automatically segmented using thresholding. A majority-voting
49 algorithm was used to generate a 'reference standard' segmentation. From the 50 CT datasets, 40 were
50 used for training and 10 for testing purposes. For cerebral structures, 34 ^{18}F -FDG PET/MRI brain image
51 volumes were used from 10 HC (5M/5F imaged twice) and 14 non-lesional epilepsy patients (7M/7F). Only
52 ^{18}F -FDG PET images were considered for training: 24/34 and 10/34 volumes were used for training and
53 testing. The dice score coefficient (DSC) was used as the primary and the average symmetric surface
54 distance (ASSD) as a secondary metric to evaluate the automated segmentation performance. **Results:** An
55 excellent overlap between the reference labels and MOOSE-derived organ segmentations was observed:
56 92% of non-cerebral tissues showed DSC values of >0.90 , while a few organs exhibited lower DSC values
57 (e.g., adrenal glands (0.72), pancreas (0.85), and bladder (0.86)). The median DSC values of brain
58 subregions derived from PET images were lower. Only 29% of the brain segments had a median DSC of
59 >0.90 , while segmentation of 60% of regions yielded a median DSC of 0.80-0.89. Results of the ASSD
60 analysis demonstrated that the average distance between the reference standard and the automatically
61 segmented tissue surfaces (organs, bones, brain regions) lies within the size of image voxels (2mm).

62 **Conclusion:** The proposed segmentation pipeline allows automatic segmentation of 120 unique tissues
63 from whole-body ^{18}F -FDG PET/CT images with high accuracy.

64

65 INTRODUCTION

66 Living organisms maintain steady internal physiological conditions through dynamic, self-regulating multi-
67 organ systemic interactions (1), also known as *homeostasis*. In healthy subjects, any notable deviation
68 from homeostasis is avoided with the aid of systemic feedback loops (2). Chronic pathologies are
69 conceived as sustained disturbances in homeostasis, which cannot be compensated by systemic
70 communications (3). Molecular imaging modalities, such as positron emission tomography (PET), can
71 provide essential insights into diverse biological processes within the human body by using highly-specific
72 radiotracers that track molecular function *in vivo* (4). Assuming that homeostasis is associated with a
73 balanced, albeit variable, glycolytic pattern, PET can help characterise bespoke feedback loops and
74 deviations that lead to pathologies. However, until recently, whole-body PET imaging protocols were
75 typically limited to only a portion of the patient's body (e.g., neck to upper thigh) due to the relatively
76 narrow axial field-of-view (FOV, 15-25cm) of PET systems. This limitation required multiple bed positions
77 to be acquired sequentially to cover the axial field of investigation. Nonetheless, this acquisition mode
78 failed to fully harness the multi-systemic physiological information provided by PET imaging (5).

79 With the recent advent of large axial FOV PET/CT systems (>70 cm) (6–8), the opportunity arose
80 to acquire total-body (TB) PET images with only 1-2 bed positions, facilitating multi-organ system analysis.
81 Such systemic analysis might allow the investigation of multi-organ interactions in various pathologies,
82 such as those associated with cancer (9), cachexia (10,11), metabolic syndrome (12), or the more recent
83 COVID-19 virus (13). However, the amount of data generated by this new generation of PET/CT systems
84 are too large to be adequately analysed without automated processing pipelines.

85 In response, we developed a multi-organ objective segmentation (MOOSE) tool, an open-source
86 software framework based on data-centric Artificial Intelligence (14) (AI) principles (Supplemental Fig. 1)
87 to allow fully-automated generation of a subject-specific total-body (TB) ¹⁸F-FDG PET/CT tissue-map
88 consisting of over 100 different tissue types. The development of such a software tool dramatically

89 increases the amount of information that can be efficiently extracted from PET data. Further, such a tool
90 provides a means to observe normal physiology and pathological conditions globally, permitting systems-
91 level investigations into human physiology. For example, when applied in a clinical setting, this approach
92 will allow physicians to automatically generate a list of standard uptake values (SUV) for all organs of
93 interest, which might provide auxiliary information during the diagnostic process. In addition, the
94 automated generation of a complete set of organ-specific SUVs lends itself well to AI-supported diagnostic
95 screening, allowing organ SUV ratios to be compared across subjects and alerting the physician about
96 potential secondary pathologies.

97

98 MATERIALS AND METHODS

99 Overall segmentation strategy

100 Our approach is based on the latest state-of-the-art nnU-Net segmentation framework (15). More
101 importantly, we propose a data-centric approach (14,16) where the network model is fixed, and the data
102 is iteratively augmented to increase the performance of the AI system. As such, the model's performance
103 is continuously monitored. As new data deviating from the training dataset's characteristics enter the
104 processing stream, the model is retrained to enhance performance.

105 Data

106 Two different types of datasets were used for the development of a software tool able to segment
107 both cerebral (83 regions) and non-cerebral structures (37 tissues).

108 For training and evaluation of non-cerebral structures, 50 whole-body low-dose CT datasets were used.
109 Among these 50 datasets, 30 CT images were acquired from healthy volunteers (14M/16F, 47 ± 13 years)
110 using the uEXPLORER total-body PET/CT system (17). The remaining 20 datasets belonged to a
111 retrospective patient cohort from a Siemens TruePoint TrueView (TPTV) PET/CT system (14M/6F, 67 ± 12
112 years). The non-cerebral tissues atlas consists of 13 abdominal organs, 20 bone segments, subcutaneous
113 fat, visceral fat, psoas, and skeletal muscle (Supplemental Table 1, Supplemental Fig. 2).

114 An expert segmentation panel comprised of four physicians and four medical students (final year)
115 was responsible for the manual segmentation of all non-cerebral structures, except for subcutaneous and
116 visceral fat and skeletal muscle, which were outlined using an established thresholding method (18). The
117 physicians were responsible for segmenting the abdominal organs and psoas while the students
118 generated the bone segments. From the 50 datasets, 40 were used for training, and 10 were used for
119 testing (hold-out dataset) purposes.

120 For training and evaluation of cerebral structures, we used 34 ^{18}F -FDG PET/MRI brain datasets (10 healthy
121 controls test-retest: 5M/5F, 27 ± 7 years and 14 non-lesional epilepsy patients: 7M/7F, 29 ± 9 years)

122 (19,20). The cerebral atlas consisted of 83 brain subregions (Supplemental Table 1), automatically created
123 from PET data in combination with T1-MR images and the Hammersmith atlas (21). In short, subject-
124 specific T1-MR images were normalized to MNI space using SPM 12 (22). The obtained (inverse) transform
125 was then used to spatially transform brain regions of the Hammersmith atlas into the individual subject's
126 native space, yielding 83 subject-specific cerebral sub-regions which were transferred to coregistered PET
127 image volumes. Of the 34 datasets, 24 and 10 were used for training and testing, respectively.

128 **Reference standard generation**

129 To address inter-variability issues of the organ segmentation, the Simultaneous Truth and
130 Performance Level Estimation (STAPLE) algorithm (23) was employed to generate reference volumes for
131 further performance assessment. Each reference volume represents a probabilistic estimate of the “true”
132 segmentation as well as a measure of multi-operator segmentation performance, (STAPLE variance). The
133 STAPLE method was not employed for reference segmentations derived using automatic (brain atlas) or
134 semi-automatic (thresholding) methods.

135 **U-Net-based semantic segmentation**

136 The nnU-Net implementation of the generic U-Net architecture is a self-configuring method for
137 deep learning-based biomedical image segmentation. This implementation exhibits strong performance
138 by retaining the original U-Net-like architecture while automating the complex process of manual
139 hyperparameter configuration (15).

140 In our implementation, training of the nnU-Net was performed separately for the following four structure
141 classes: (i) 13 abdominal organs and psoas, (ii) 20 bone structures, (iii) 83 brain regions (iv) fat
142 (subcutaneous and visceral) and skeletal muscle. Segmentation of non-cerebral tissues was performed
143 based on CT data, whereas segmentation of cerebral regions was carried out using ¹⁸F-FDG PET images.

144 **Assessment of deviation from training dataset distribution**

145 It is unlikely that any training dataset will be sufficient to fully capture the variability encountered
146 in clinical routine. Accordingly, a data-centric approach is necessary, permitting continuous monitoring of
147 segmentation performance so that data that substantially deviates from the original training data
148 distribution (i.e., Out-of-distribution (OOD) data) is detected. Erroneous segmentation results obtained
149 for such data will then require manual correction by a human expert. Once corrected, this data can be
150 appended in suitable quantities to the initial training dataset for retraining purposes.

151 Since continuous operator-based monitoring of segmentation performance is untenable in clinical
152 routine, we developed an automated error analysis routine that detects OOD datasets based on
153 morphometric analysis of organ shapes (e.g., elongation, volume, area, maximum and minimum bounding
154 box diameter), which were determined for each STAPLE-derived segmentation of structures, and a
155 normative morphological feature database was generated. Upon segmentation of a new dataset,
156 morphological features for each segmented structure are calculated and compared to the normative
157 morphology database, yielding a distance (Z-score) in “similarity space” for each structure. The Z-score
158 reflects the difference between the shapes of the segmented structure in comparison to its normative
159 value obtained from the training datasets. In our implementation, we chose a value of $Z = 1.5$ as the cut-
160 off for OOD labelling.

161 **Algorithm performance vs. training sample size**

162 Primary performance assessment of the automated segmentation (MOOSE) was performed for
163 all structures using the Dice score coefficient (DSC) (24). A DSC value of 1.0 with respect to STAPLE
164 indicates perfect overlap and 0 indicates no overlap. In addition, the average symmetric surface distance
165 (ASSD) (25) was used as a secondary metric, representing the average distance (in mm) between surface
166 voxels of the standard and the automated segmentation.

167 To assess the segmentation performance as a function of training sample size, we calculated for
168 each non-cerebral structure the DSC and the ASSD using the segmented volumes derived using 10 (D10),

169 20 (D20), and 40 (D40) training data sets, respectively. A similar analysis was performed for cerebral
170 regions with 8 (D'8), 16 (D'16), and 24 (D'24) data sets. In both instances, cases were randomly selected
171 from the whole datasets (50 cases for non-cerebral and 34 cases for cerebral structures). The testing (hold-
172 out) dataset included 10 cases that were not part of the training sets in both instances.

173 **Algorithm performance vs. training dataset variability**

174 To investigate the effect of training dataset variability on segmentation performance, we
175 performed a series of training/test runs using various mixtures of two datasets that differed significantly
176 with respect to arm position (either arms down or crossed on chest, Supplemental Fig. 3). We created
177 four subsets of training datasets, each with a total sample size of 20. The first dataset consisted of 20 low-
178 dose CT images with arms down (SMS20). The other three training datasets included mixtures of images:
179 (i) MIX2-18 (18 arms down, 2 crossed), MIX5-15 (15 arms down, 5 crossed), and MIX10-10 (10 arms down,
180 10 crossed). Networks trained on these four training datasets were then used to segment 10 test datasets
181 that included only images with crossed arm positions (X10). The following four (training test) scenarios
182 were investigated: [SMS20→X10], [MIX2-18→X10], [MIX5-15→X10] and [MIX10-10→X10]. Segmentation
183 results were assessed separately for bone structures of the arm (radius, ulna, carpal, metacarpal, fingers)
184 and for all other bone structures (that did not differ positionally). This analysis provided information
185 regarding the necessary variability in the training dataset required to segment OOD data accurately.

186 **Algorithm performance for clinical OOD datasets**

187 We applied the trained network to two small pathological cohorts that were not part of the initial
188 training set: three lymphoma and three mesothelioma lung cancer cases. The intent was to assess the
189 performance of MOOSE on clinical datasets that differ significantly from the training data distribution.
190 The segmentation quality was evaluated based on similarity space analysis (Z-scores). OOD datasets with
191 incorrect segmentations were manually corrected, and the corrected segmentations were then appended
192 to the original training datasets for retraining purposes.

193 **Statistical assessment**

194 A paired t-test was applied to determine whether DSC values differed significantly between the
195 various training sample sizes and to investigate the effect of training dataset variability (either fully OOD
196 or mixed) on DSC values. In addition, a correlation analysis (Pearson’s rho) was performed to investigate
197 the relationship between STAPLE variance and the DSC values associated with the best (D40) training
198 sample size. A similar analysis was also performed using the ASSD metric.

199 **Software tool implementation**

200 Our processing pipeline is based on Python, and C++ programming languages, with the nnU-Net
201 framework representing the segmentation backbone, built using PyTorch 1.6.031 (26). Similarity space
202 was implemented using the morphometric capabilities of SimpleITK 2.1.0 (27), and manual cleaning of
203 erroneous segmentation results was performed using 3D Slicer (28) (Version 4.11.20210226).

204

205 RESULTS

206 Effect of training data size on segmentation performance

207 The majority of non-cerebral tissues (81%) were segmented with high accuracy (DSC>0.90) as seen
208 from Fig. 1. Of note, DSC analysis generally showed an excellent overlap between STAPLE-derived
209 reference and organ segmentations based on 10 (D10), 20 (D20), and 40 (D40) training data sets. This
210 excellent overlap was confirmed based on ASSD analysis, yielding average distances of 1.40 ± 1.29 , $1.05 \pm$
211 1.26 and 0.68 ± 0.52 mm for D10, D20 and D40, respectively. However, the performance of the automated
212 segmentation was suboptimal for a small group of organs (highlighted in red, Fig. 1), with low median DSC
213 and high SD values (Supplemental Fig. 4A), such as the adrenal glands (DSC=0.72), pancreas (0.85), and
214 bladder (0.86). Subsequent correlation analysis of the STAPLE variance and the DSC values derived from
215 the (best) D40 training set is shown in Supplemental Fig. 4B. The graph indicates an overall highly
216 significant negative correlation ($\rho = -0.79$, $p = 0.002$) with the three identified regions showing high
217 STAPLE variance. This significant correlation with the STAPLE variance was also reproduced using the ASSD
218 metric ($\rho = 0.60$, $p = 0.042$, Supplemental Fig. 5), indicating that accurate segmentation of this subset
219 of regions is challenging even for human experts.

220 The segmentation performance for bone structures was similar to that of the abdominal organs
221 (Fig. 2). Again, one notes an excellent overlap between the reference structure volumes and those
222 obtained using the automated segmentation based on 10 (D10), 20 (D20), and 40 (D40) training data sets
223 (ASSD of 1.63 ± 3.01 , 1.61 ± 3.14 and 0.83 ± 0.76 mm, respectively), except for a small number of bone
224 structures with either low DSC mean or high SD values (Supplemental Fig. 6). These structures were the
225 carpal bones, metacarpal bones, and phalanges of the toes. Removal of these organs resulted in a similar
226 segmentation performance between D20 and D40 ($p = 0.07$), with segmentation based on D10 remaining
227 significantly worse than D20 ($p = 0.016$) and D40 ($p = 0.010$).

228 Although median DSC values of brain subregions derived from PET images were relatively low
229 (only 29% of brain segments had median DSC values >0.90, see Supplemental Fig. 7), ASSD values showed
230 sub-voxel differences between the template regions and the automated segmentation, with similar
231 performance across the D10 ($0.52 \pm 0.35\text{mm}$), D20 ($0.53 \pm 0.41\text{mm}$) and D40 ($0.46 \pm 0.27\text{mm}$) datasets.

232 **Effect of training dataset variability on segmentation performance**

233 Results of dataset variability analysis are shown in Fig. 3. The figure indicates that segmentation
234 of structures that substantially deviate from the ‘expected’ position in the training datasets was
235 suboptimal. However, by including at least two cases that match the deviant position to the training
236 dataset resulted in a significant performance improvement. Specifically, DSC values for bones of the arm
237 were significantly lower for the fully OOD [SMS20→X10] scenario (DSC = 0.87 ± 0.12) as compared to the
238 three scenarios that included 10% [MIX2-18→X10] (DSC = 0.92 ± 0.06 ; $p = 0.04$), 25% [MIX5-15→X10] (DSC
239 = 0.940 ± 0.003 ; $p = 0.01$) and 50% [MIX10-10→X10] (DSC = 0.91 ± 0.04 ; $p = 0.04$) of cases that matched
240 the deviant position. In addition, the COV for DSC values derived from the three mixed training datasets
241 was significantly lower (COV = 6.6%, $p = 0.01$; COV = 3.3%, $p = 0.03$; COV = 4.3%, $p = 0.01$) as compared to
242 the COV for DSC values derived using the fully OOD training dataset (COV = 13.5%). In comparison,
243 performance of all four scenarios for bone structures that were matched in position between the training
244 and test datasets was similar, with DSC values of >0.95 (Fig. 3).

245 **Detection of OOD segmentation errors**

246 Application of similarity space analysis identified segmentation errors in clinical datasets that
247 included various anatomical pathologies, representing OOD datasets for specific organs. This was clearly
248 demonstrated by applying the initially trained neural network to two distinct OOD datasets (lymphoma
249 and mesothelioma) that were not part of the initial training set. Specifically, all lymphoma patients
250 presented with splenomegaly which led to its incorrect classification as a “liver” and “spleen” (Fig. 4A).
251 Following manual correction (time required ~3min per case), we appended two corrected data sets to the

252 original training set to retrain the neural network. The retrained neural network correctly segmented the
253 abnormally enlarged spleen in the third lymphoma patient which was used as a hold-out dataset (Fig. 5).

254 Similarly, the large tumor mass in the lungs of mesothelioma patients was incorrectly classified as
255 part of the liver, heart, and bladder (Fig. 4B). Again, similarity space analysis identified the incorrect
256 segmentation and labelled the dataset as representing an OOD image pattern (Fig. 6A). Following manual
257 correction of two out of three patients, these two cases were again appended to the training data set,
258 and the neural network was retrained using the extended training set. Once again, we determined an
259 improvement in the segmentation performance of the third (uncorrected) data set (Fig. 6).

260

261 **DISCUSSION**

262 Hybrid molecular imaging modalities, such as ^{18}F -FDG PET/CT allow the investigation of multi-
263 organ systemic interactions through which living organisms maintain homeostasis and allostasis. The
264 resulting images are not mere pictures - they represent a rich pallet of multi-dimensional data (29). By
265 systemically parcellating these datasets into respective organ/tissue classes, one can, in theory, study
266 system-level interactions in detail between the various homeostatic and allostatic networks, allowing a
267 better understanding of pathological abnormalities in vivo. Nevertheless, manual segmentation of various
268 tissues in the human body is not tenable, either in research applications or in clinical routine.

269 To bridge this gap, we developed a fully-automated segmentation pipeline (MOOSE) that allows
270 the creation of subject-specific multi-tissue FDG PET/CT atlases (Supplemental Fig. 2). These tissue-maps
271 enables the extraction of subject-specific functional information from molecular imaging data with
272 minimal additional effort for further analysis. We based the developed segmentation pipeline on the
273 latest state-of-art nnU-Net architecture (15) and demonstrated that robust training of the convolutional
274 neural network could be achieved with as low as 20 datasets, provided that sufficient variability in the
275 training dataset is present. In addition, our results support the concept of data-centric AI, which focuses
276 primarily on data quality rather than quantity.

277 In general, MOOSE performed reasonably well in segmenting most of the non-cerebral tissues
278 while exhibiting poorer segmentation performance on selected organs, such as thyroid, adrenal gland and
279 bladder. Our correlation analysis revealed a significant negative correlation between the STAPLE variance
280 and the DSC values derived from the (best) D40 training set (Supplemental Fig. 4B). This result suggests
281 that, due to a combination of small organ size, low contrast and increased noise levels present in low-
282 dose CT images, accurate segmentation of bespoke structures is challenging even for human experts.

283 **AI, PET imaging, and systems biology**

284 The ultimate objective of the developed multi-organ and tissue segmentation methodology is to
285 promote the concept of whole-person research (30) and systems biomedicine (31) through whole/total-
286 body ¹⁸F-FDG PET/CT imaging. With the advent of large axial FOV PET/CT systems, most or all organs can
287 be simultaneously imaged, therefore allowing an improved evaluation of interactions between organs in
288 both healthy and diseased states. We envision that through automated extraction of rich physiological
289 information inherent in PET/CT data (e.g., organ SUVs), disease-specific “metabolic fingerprints” can be
290 derived that uniquely characterize diverse pathologies affecting system-level organ interaction
291 (Supplemental Fig. 8). Such an analysis might uncover metabolic dependencies among sets of organs and
292 might provide novel insights into metabolic pathway dysregulation associated with disease progression.
293 Moreover, given the fact that non-cerebral tissues are segmented directly from CT data, this technique is
294 insensitive to variations in PET tracer uptake patterns, thus allowing the study of diverse system-level
295 functional processes using a multitude of function-specific radiotracers.

296 **Training of neural networks using sparse datasets**

297 It is commonly assumed that the performance of a neural network increases with the size of the
298 training set. Therefore, most non-healthcare image classification applications are trained on vast numbers
299 of training cases (e.g., ImageNet (32)). However, creating large training datasets in the medical field is
300 problematic, as manual curation of medical images is highly time-consuming and heavily dependent on
301 domain-specific human expert knowledge. In this study, 50 medical image datasets were manually
302 segmented (into 120 objects for each dataset) by medical professionals. This process required significant
303 personal effort by each expert and took several months to complete. Such an effort cannot be expected
304 to be repeated numerous times when additional silos of data (possibly with a different distribution)
305 become available.

306 In recognition of this methodological constraint, we investigated the effect of both training
307 sample size and training sample variability on segmentation performance. Our results demonstrate that

308 segmentation performance is primarily dependent on whether the training dataset allows the correct
309 identification of several distinct unique image patterns, each characterized by a mean spatial pattern and
310 the associated variance (**Fig 4-6**). This insight also explains why more cases are usually preferred, as it is
311 likely that a greater number of unique image patterns can be captured using a larger dataset. However,
312 the number of images needed per unique pattern is not evident. Our results suggest that accurate
313 segmentation of abnormal image patterns is viable, provided that the training data includes a small
314 number (2-4) of cases that establish a distinct image pattern with the associated morphological variance.

315 **A data-centric approach to segmentation**

316 Over the long run, any clinically viable medical image segmentation method will require a system
317 where incoming data is constantly utilized to adjust model parameters to accommodate changing data
318 distributions. To accomplish this, the implemented data-centric approach executes two operations: first,
319 it performs active monitoring of segmentation performance followed by the users input to correct the
320 erroneous segmentation, and second, implements periodic updating of model parameters through
321 retraining of the neural network using an updated training set (which includes the manually corrected
322 OOD data).

323 In particular, segmentation performance is continuously monitored in similarity space, and
324 feedback regarding segmentation accuracy is provided to the physician in the form of tissue-specific Z-
325 scores that signal potential deviations from tissue shape/position in the normative training data
326 distribution. Based on this analysis, all tissues that are judged to be out-of-distribution ($Z > 1.5$) are
327 flagged, and the physician is prompted for corrective action. This approach ensures adequate
328 segmentation of all tissues present in abnormal datasets and provides important curated data for future
329 retraining of the neural network. Moreover, this strategy addresses potential segmentation problems
330 right when they occur in the processing pipeline, when corrective actions can be carried out most
331 efficiently and with the least effort.

332 It is important to note that the presented segmentation framework bears its challenges. First and
333 foremost, this methodology mandates a high-performance workstation which might be cost-prohibitive.
334 Our network training was performed on a dedicated server (Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz
335 32 CPU cores, 256 GB of RAM, and 1 x Nvidia GeForce RTX 3090 Ti), allowing the generation of one TB
336 ¹⁸F-DG PET/CT tissue-map from an individual whole-body PET/CT dataset in about 30 minutes. Moreover,
337 once OOD data sets are collected, the neural network needs to be retrained, which took approximately
338 two days to complete using the above server configuration. Finally, there is some unavoidable subjectivity
339 in identifying OOD data sets as the cut-off defining OOD data is based on heuristics.

340

341 **CONCLUSION**

342 We present here a fully-automated, data-centric segmentation pipeline for the creation of a total-
343 body ¹⁸F-FDG PET/CT tissue-map. The generated map is modular and consists of 120 tissues and bone
344 structures, enabling the automated extraction of image information for both cerebral and non-cerebral
345 regions, potentially providing added information about secondary abnormalities during the diagnostic
346 process.

347

348 **CODE-AVAILABILITY**

349 We have named our software pipeline MOOSE¹²⁰, which stands for Multi-Organ Objective
350 Segmentation. MOOSE¹²⁰ is free and is an open-source software. All codes related to MOOSE¹²⁰ are
351 available online. All models for our application are publicly available, and a complete description of the
352 processing pipeline is available on our GitHub page (<https://github.com/QIMP-Team/MOOSE>).

353

354 **CONFLICT OF INTEREST**

355 No potential conflicts of interest relevant to this article exist.

356

357 **ACKNOWLEDGMENTS**

358 We thank Jakub Gawrylkowicz and Sebastian Gutschmayer for their support in packaging the
359 software. We extend our gratitude to Clemens Spielvogel for helpful discussions throughout the project.
360 We thank Kilian Kluge for organising the mesothelioma datasets, and finally, we thank Irene Buvat for
361 providing the lymphoma datasets to verify our hypothesis. We extend our gratitude to IBM, as this
362 research is supported through the IBM University Cloud Award.

363

364 **KEY POINTS:**

365 **QUESTION:** How can we efficiently extract diagnostic information from whole-body ^{18}F -FDG-PET/CT data?

366 **PERTINENT FINDINGS:** An automated approach for multi-organ segmentation of whole-body ^{18}F -FDG-
367 PET data is presented. It builds on the nnU-Net methodology driven by data-centric principles and
368 supports accurate segmentation of 37 extra-cerebral and 83 cerebral regions. Over 92% of the non-
369 cerebral tissues were segmented with a Dice score value of more than 0.90, while 89% of the cerebral
370 areas had a DSC of more than 0.80.

371 **IMPLICATIONS FOR PATIENT CARE:** The developed software tool increases the amount of information
372 extracted from standard, whole-body PET/CT datasets and provides means to perform system-level
373 investigations into human physiology.

374 **REFERENCES**

- 375 1. Cannon WB. The wisdom of the body. *Am J Med Sci.* 1932;184:864.
- 376 2. Goodman L. Regulation and control in physiological systems: 1960-1980. *Ann Biomed Eng.*
377 1980;8:281-290.
- 378 3. Billman GE. Homeostasis: The Underappreciated and Far Too Often Ignored Central Organizing
379 Principle of Physiology. *Front Physiol.* 2020;11:200.
- 380 4. Lammertsma AA. Forward to the Past: The Case for Quantitative PET Imaging. *J Nucl Med.*
381 2017;58:1019-1024.
- 382 5. Cherry SR, Badawi RD, Karp JS, Moses WW, Price P, Jones T. Total-body imaging: Transforming the
383 role of positron emission tomography. *Sci Transl Med.* 2017;9:eaaf6169.
- 384 6. Karp JS, Viswanath V, Geagan MJ, et al. PennPET Explorer: Design and Preliminary Performance of a
385 Whole-Body Imager. *J Nucl Med.* 2020;61:136-143.
- 386 7. Spencer BA, Berg E, Schmall JP, et al. Performance evaluation of the uEXPLORER total-body PET/CT
387 scanner based on NEMA NU 2-2018 with additional tests to characterize PET scanners with a long axial
388 field of view. *J Nucl Med.* 2021;62:861-870.
- 389 8. Prenosil GA, Sari H, Fürstner M, et al. Performance Characteristics of the Biograph Vision Quadra
390 PET/CT System with a Long Axial Field of View Using the NEMA NU 2-2018 Standard. *J Nucl Med.*
391 2022;63:476-484.
- 392 9. Zhu L, Finkelstein D, Gao C, et al. Multi-organ Mapping of Cancer Risk. *Cell.* 2016;166:1132-1146.e7.
- 393 10. Penet M-F, Winnard PT Jr, Jacobs MA, Bhujwalla ZM. Understanding cancer-induced cachexia:
394 imaging the flame and its fuel. *Curr Opin Support Palliat Care.* 2011;5:327-333.
- 395 11. Argilés JM, Busquets S, Stemmler B, López-Soriano FJ. Cancer cachexia: understanding the molecular

- 396 basis. *Nat Rev Cancer*. 2014;14:754-762.
- 397 12. Priest C, Tontonoz P. Inter-organ cross-talk in metabolic syndrome. *Nat Metab*. 2019;1:1177-1188.
- 398 13. Gupta A, Madhavan MV, Sehgal K, et al. Extrapulmonary manifestations of COVID-19. *Nat Med*.
399 2020;26:1017-1032.
- 400 14. Wu A. A Chat with Andrew on MLOps: From Model-Centric to Data-Centric AI. 2021.
- 401 15. Isensee F, Jaeger PF, Kohl SAA, Petersen J, Maier-Hein KH. nnU-Net: a self-configuring method for
402 deep learning-based biomedical image segmentation. *Nat Methods*. 2021;18:203-211.
- 403 16. Motamedi M, Sakharykh N, Kaldewey T. A Data-Centric Approach for Training Deep Neural
404 Networks with Less Data. *arXiv [csAI]*. October 2021.
- 405 17. Badawi RD, Shi H, Hu P, et al. First Human Imaging Studies with the EXPLORER Total-Body PET
406 Scanner. *J Nucl Med*. 2019;60:299-303.
- 407 18. Weston AD, Korfiatis P, Kline TL, et al. Automated Abdominal Segmentation of CT Scans for Body
408 Composition Analysis Using Deep Learning. *Radiology*. 2019;290:669-679.
- 409 19. Traub-Weidinger T, Muzik O, Sundar LKS, et al. Utility of Absolute Quantification in Non-lesional
410 Extratemporal Lobe Epilepsy Using FDG PET/MR Imaging. *Front Neurol*. 2020;11:54.
- 411 20. Sundar LKS, Muzik O, Rischka L, Hahn A. Promise of Fully Integrated PET/MRI: Noninvasive Clinical
412 Quantification of Cerebral Glucose Metabolism. *Journal of Nuclear*. 2020.
- 413 21. Hammers A, Allom R, Koepp MJ, et al. Three-dimensional maximum probability atlas of the human
414 brain, with particular reference to the temporal lobe. *Hum Brain Mapp*. 2003;19:224-247.
- 415 22. Ashburner J. SPM: a history. *Neuroimage*. 2012;62:791-800.
- 416 23. Warfield SK, Zou KH, Wells WM. Simultaneous truth and performance level estimation (STAPLE): an

417 algorithm for the validation of image segmentation. *IEEE Trans Med Imaging*. 2004;23:903-921.

418 24. Dice LR. Measures of the Amount of Ecologic Association Between Species. *Ecology*. 1945;26:297-
419 302.

420 25. Yeghiazaryan V, Voiculescu I. Family of boundary overlap metrics for the evaluation of medical image
421 segmentation. *J Med Imaging (Bellingham)*. 2018;5:015006.

422 26. Paszke A, Gross S, Massa F, et al. PyTorch: An Imperative Style, High-Performance Deep Learning
423 Library. *arXiv [csLG]*. December 2019.

424 27. Lowekamp BC, Chen DT, Ibáñez L, Blezek D. The Design of SimpleITK. *Front Neuroinform*. 2013;7:45.

425 28. Fedorov A, Beichel R, Kalpathy-Cramer J, et al. 3D Slicer as an image computing platform for the
426 Quantitative Imaging Network. *Magn Reson Imaging*. 2012;30:1323-1341.

427 29. Gillies RJ, Kinahan PE, Hricak H. Radiomics: Images Are More than Pictures, They Are Data. *Radiology*.
428 2016;278:563-577.

429 30. NCCIH strategic plan FY 2021–2025. NCCIH. [https://www.nccih.nih.gov/about/nccih-strategic-plan-](https://www.nccih.nih.gov/about/nccih-strategic-plan-2021-2025)
430 2021-2025.

431 31. Hacker M, Hicks RJ, Beyer T. Applied Systems Biology—embracing molecular imaging for systemic
432 medicine. *Eur J Nucl Med Mol Imaging*. 2020;47:2721-2725.

433 32. Deng J, Dong W, Socher R, Li L-J, Li K, Fei-Fei L. ImageNet: A large-scale hierarchical image database.
434 In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. ; 2009:248-255.

435

436 **ETHICS DECLARATION**

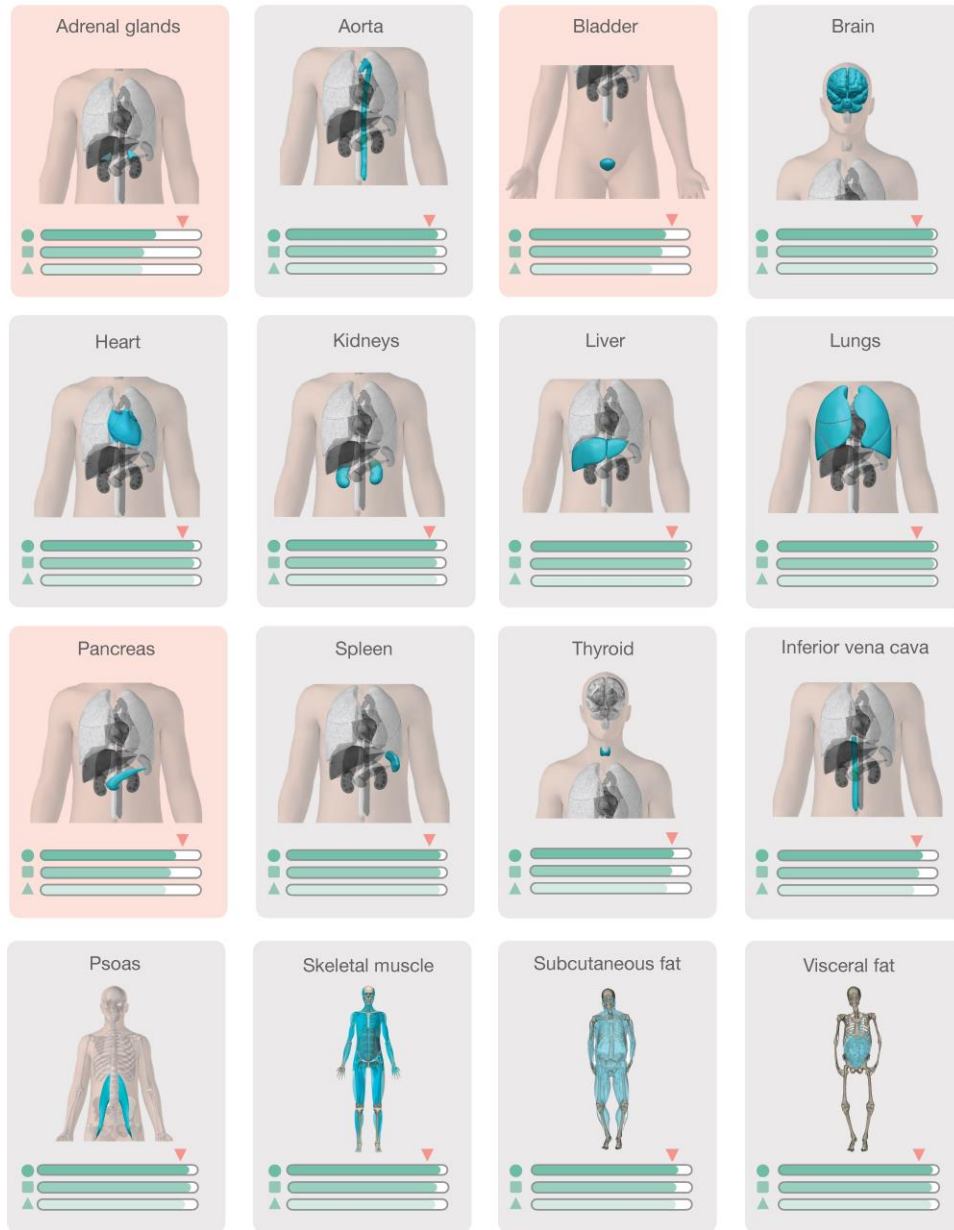
437 All the data utilized in this study were acquired in accordance with the Declaration of Helsinki. Written
 438 informed consent was obtained from all the subjects prior to the examinations.

Dataset	Acquisition system	Institutional Review Board	Reference number
34 ¹⁸ F-FDG PET/MR brain datasets	Siemens Biograph mMR, Siemens Healthineers	Medical University of Vienna	EK1960/2014
30 low-dose healthy control CT datasets	uEXPLORER, United Imaging Healthcare	The University of California at Davis	I1341792-18
20 low-dose mixed pathological Siemens CT datasets	Siemens Biograph mCT TruePoint TrueV, Siemens Healthineers	Medical University of Vienna	EK1649/2016
Three Lymphoma datasets	Philips Gemini GXL16, Philips Medical Systems	Protection des Personnes Sud-Est III, Hôpital HOTEL-DIEU, Place de l'Hôpital	Etude REMARC Réf: 2009 - 006B; Eudract N°: 2008-008202-52.
Three mesothelioma datasets	Siemens Biograph mCT TruePoint TrueV, Siemens Healthineers	Medical University of Vienna	EK1649/2016

439

440 **FIGURES**

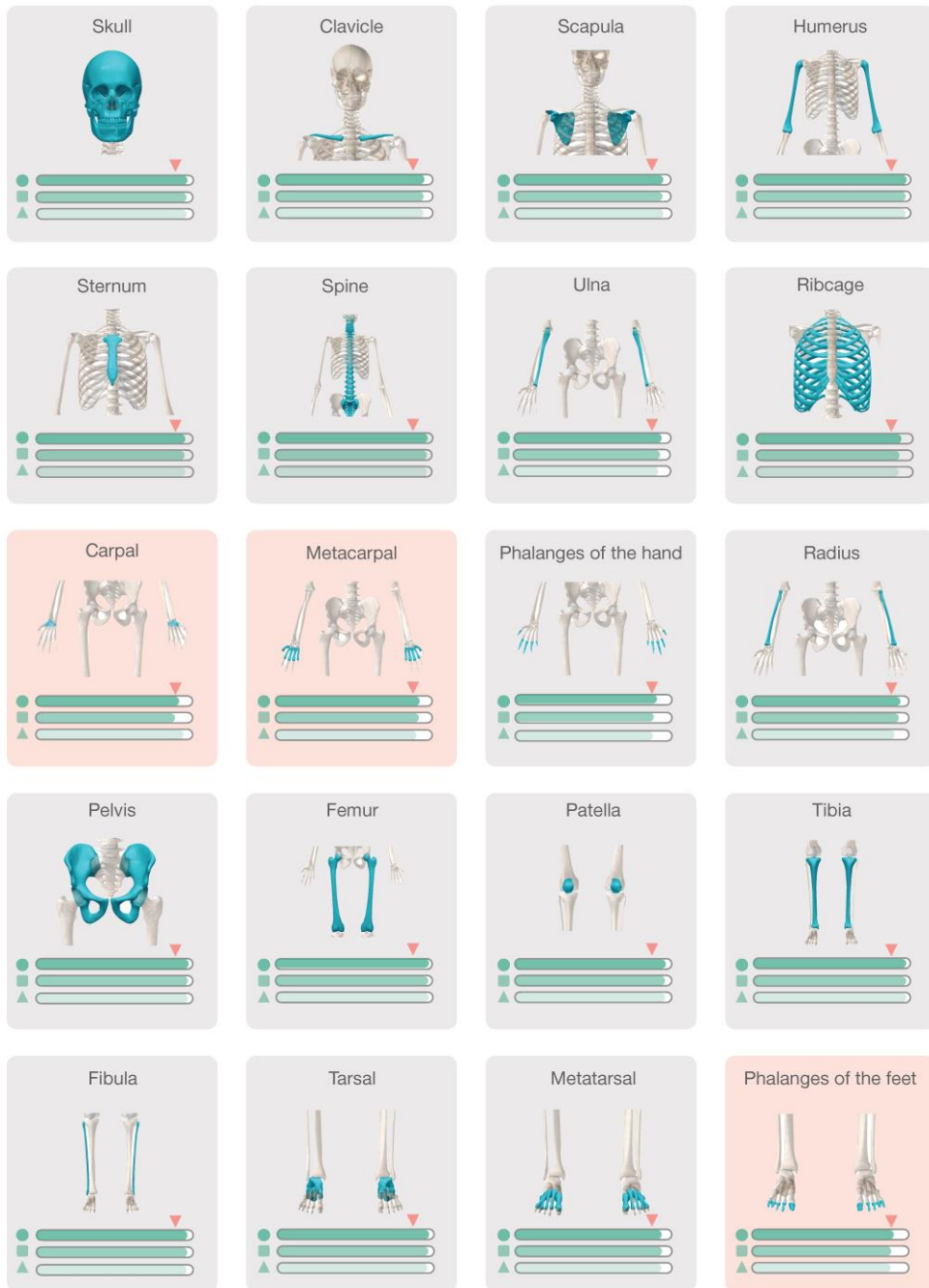
Dice-score performance for non-cerebral organs as a function of training dataset size



441

442 **Figure 1.** Median DSCs of abdominal organs (10 test datasets) were obtained from models based on three
 443 separate training subsets: D40 (circle), D20 (square), D10 (triangle). The inverted triangle (pink) indicates
 444 the 0.90 mark. Red background highlights organs characterised by low median DSCs (<0.90) and high
 445 standard deviation (SD, see **Supplemental Fig. 4**).

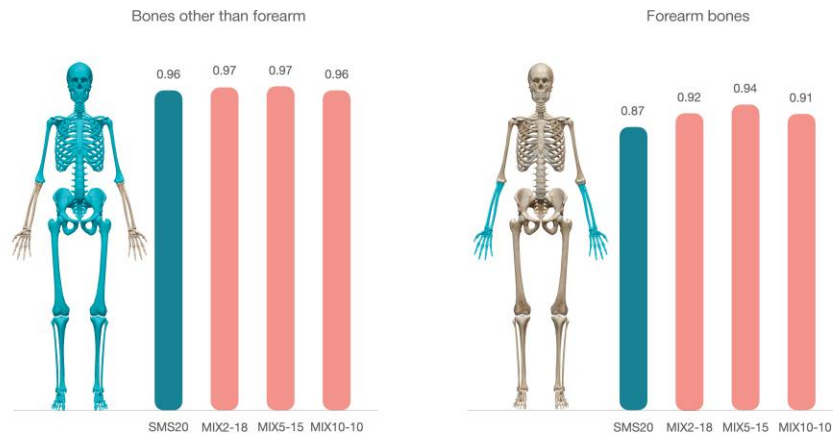
Dice-score performance for bones as a function of training dataset size



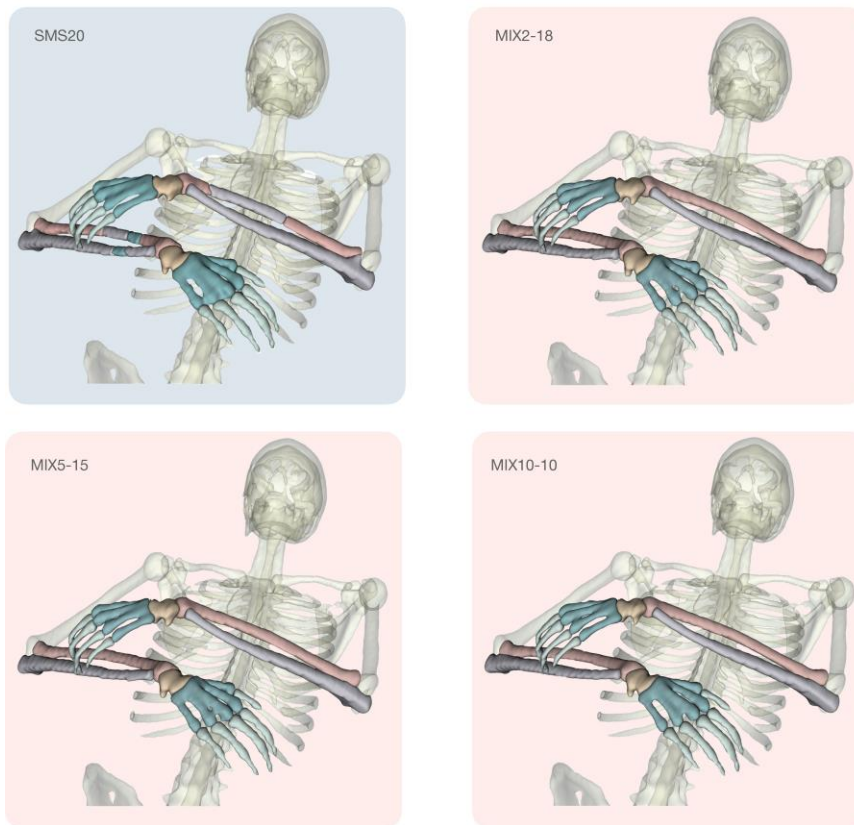
446

447 **Figure 2.** Median DSCs of bone structures (10 test datasets) as obtained from models based on three
 448 separate training subsets: D40 (circle), D20 (square), and D10 (triangle). The inverted triangular marker
 449 indicates the 0.90 mark. Red background highlights bones characterised by low median DSCs (<0.90).
 450

451

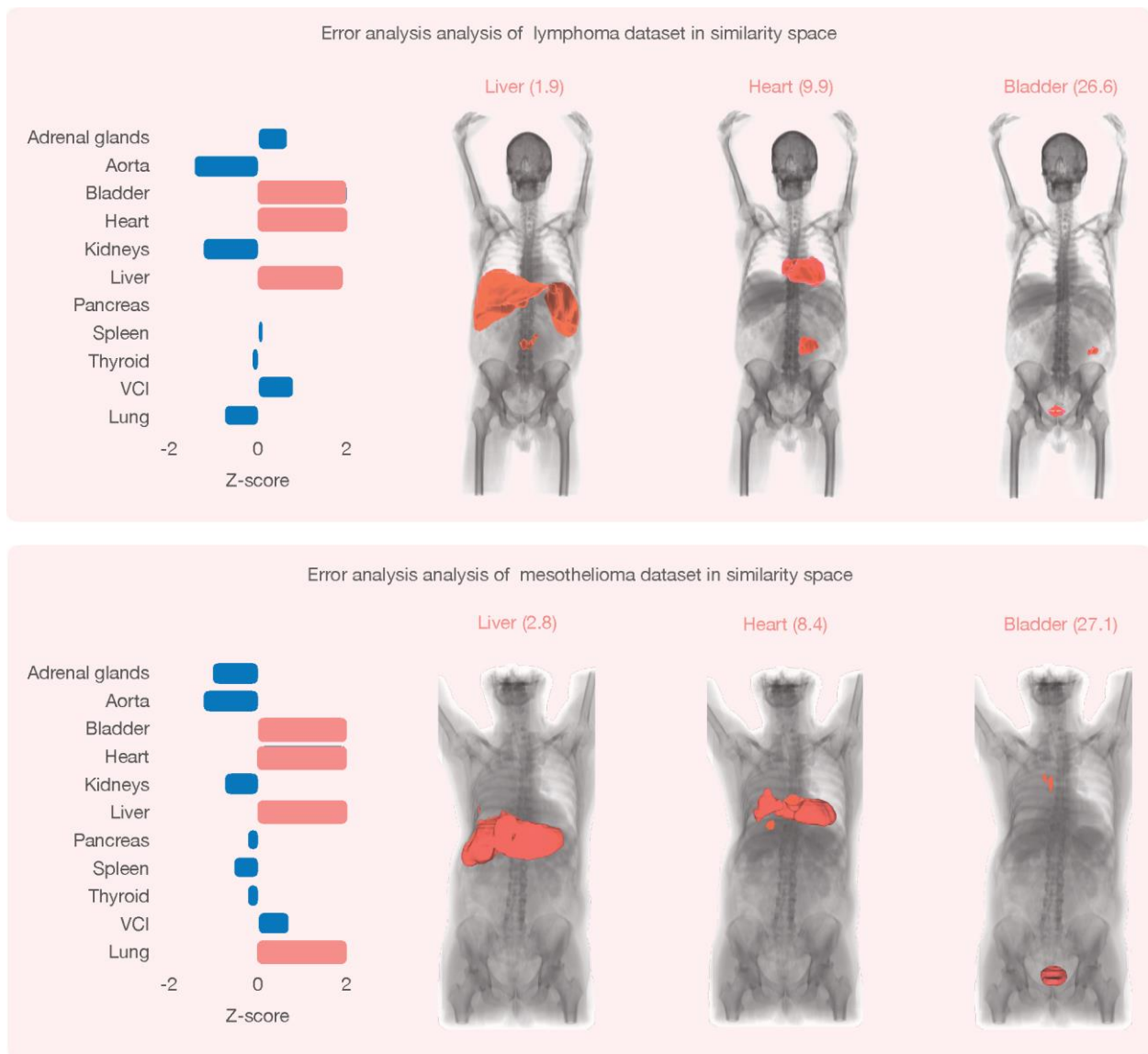


Forearm bone analysis for an individual subject



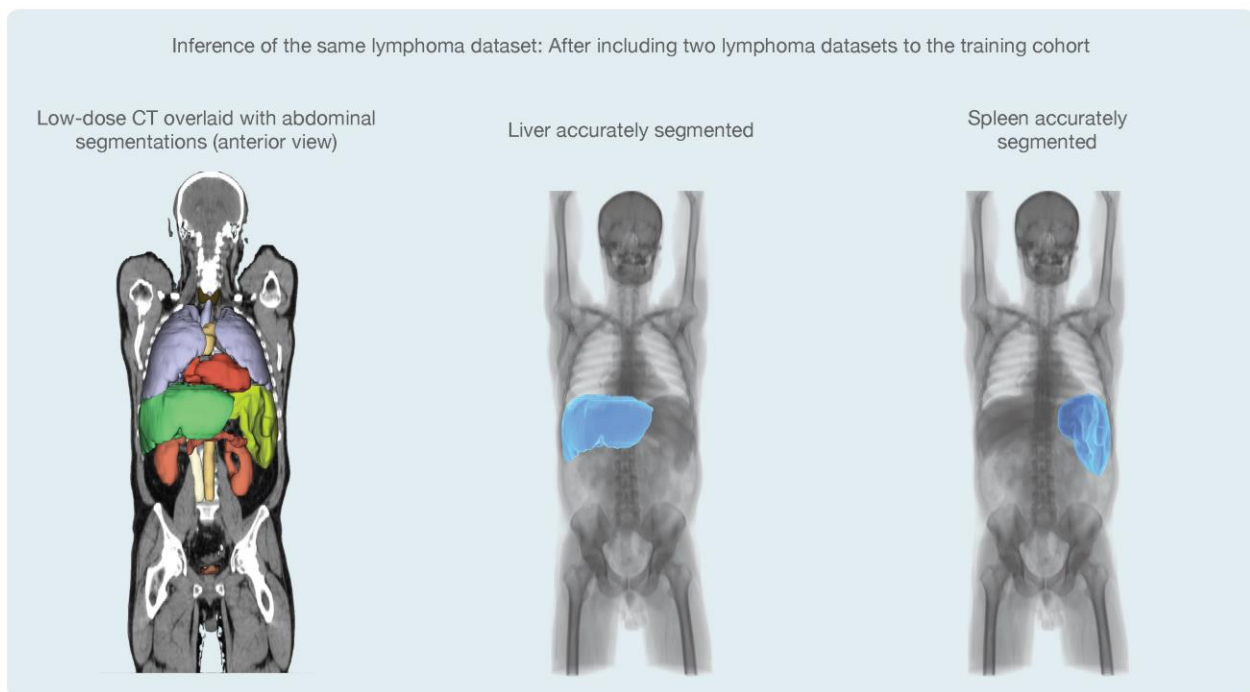
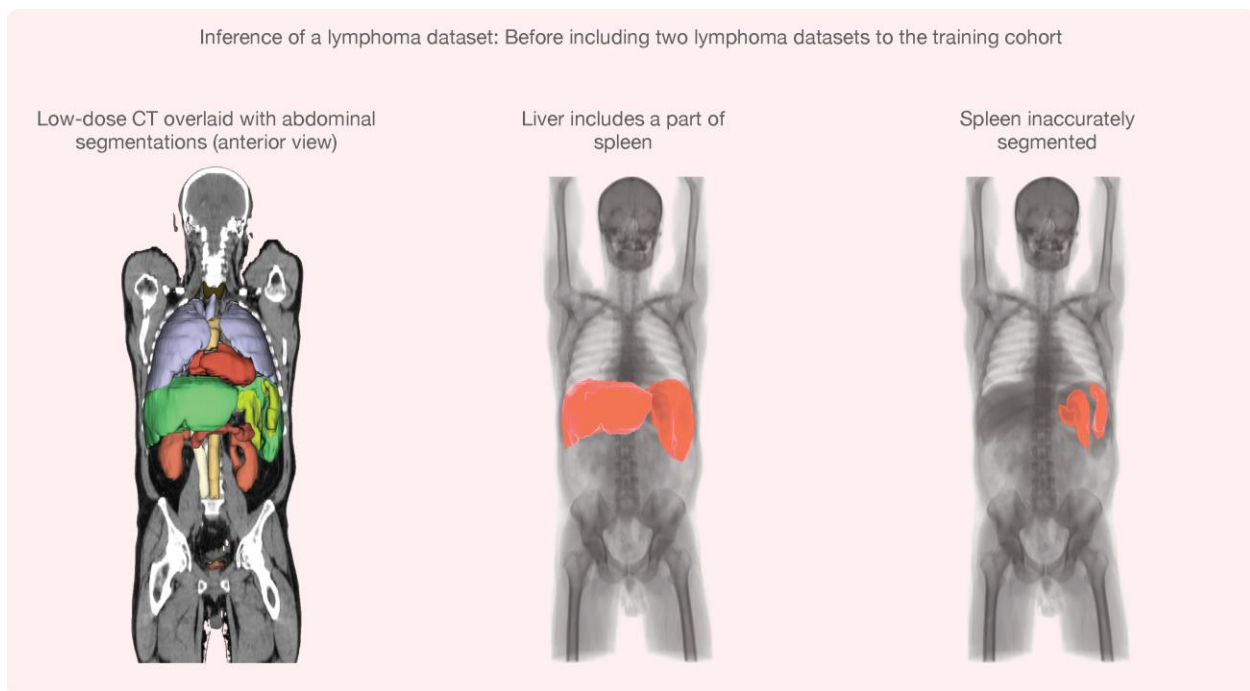
452

453 **Figure 3.** (A) Bar graph demonstrating a similar performance of different models for bone segmentation
 454 other than forearm bones. Green bar depicts the homogenous training dataset (SMS-20), whereas mixed
 455 training datasets (MIX2-18, MIX5-15, MIX10-10) are represented by red bars. (B) Bar graph showing
 456 segmentation performance of forearm bones. A significant performance improvement is seen in the
 457 mixed training datasets (red bars) compared to the homogeneous training dataset (green bar). (C)
 458 Forearm bone analysis of an individual subject. The images demonstrate that the forearm bones are
 459 incorrectly segmented in the case of the SMS20 (green background) model, whereas all mixed models
 460 accurately segmented the forearm bones (red background).



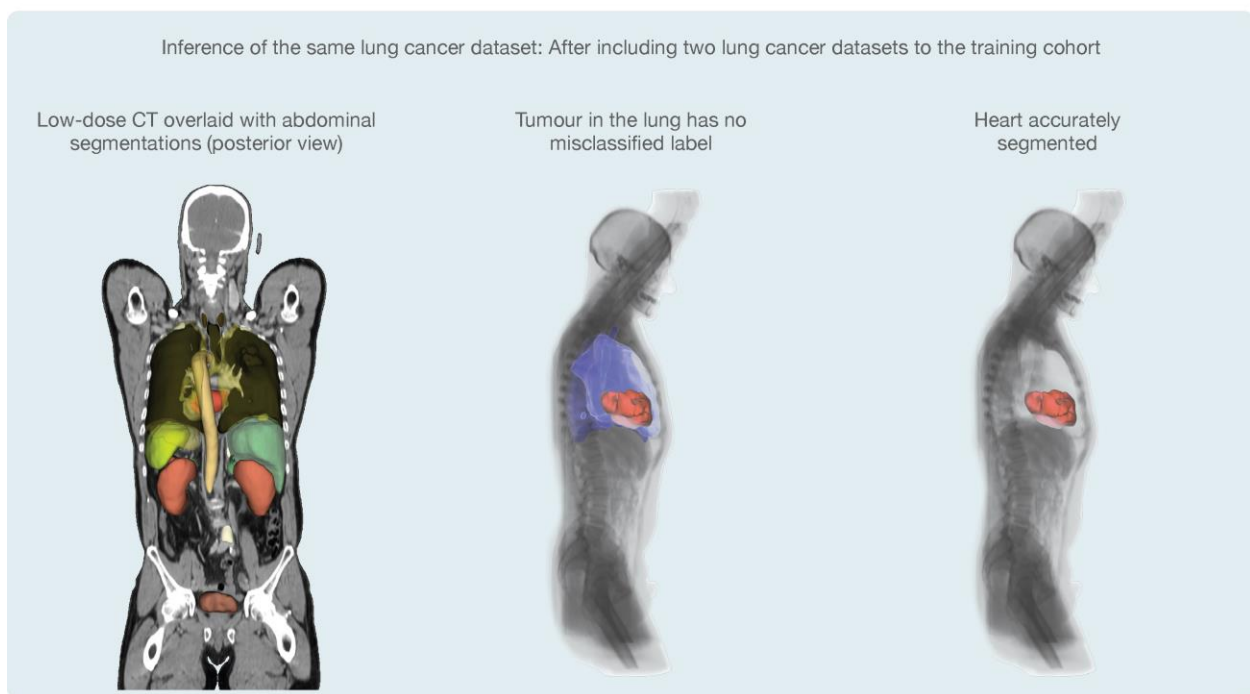
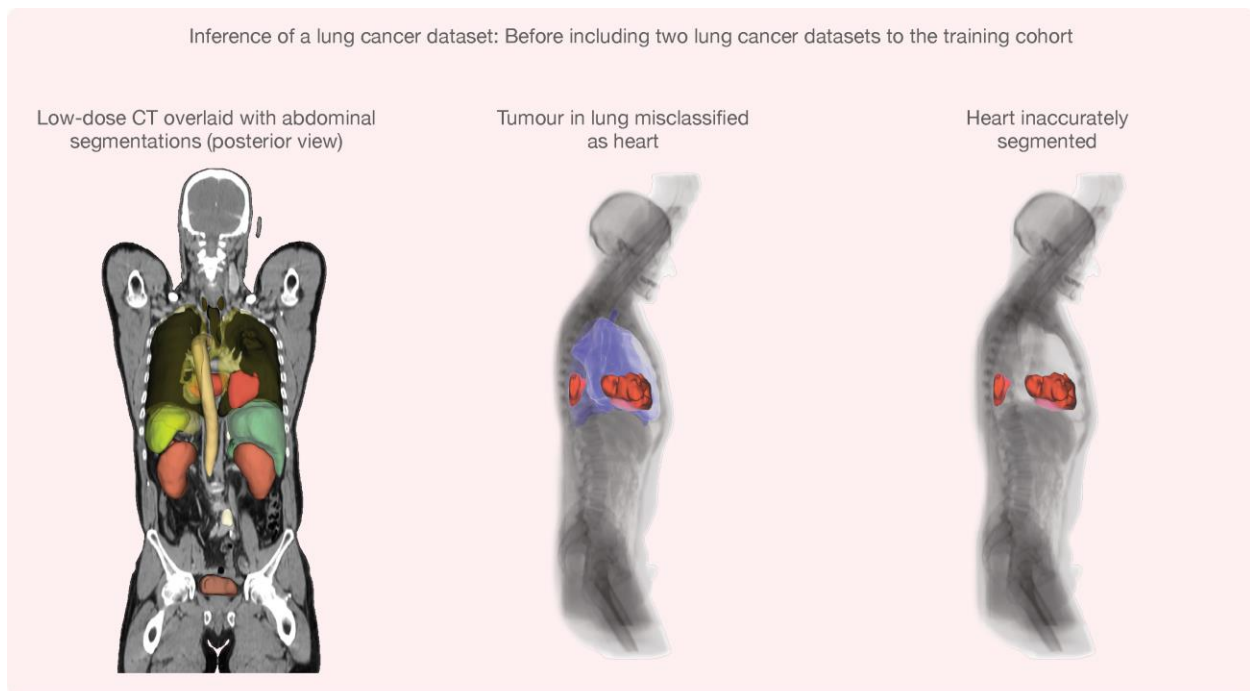
461

462 **Figure 4** (A) Error analysis in similarity space for a representative lymphoma patient. Horizontal bars depict
 463 distance in similarity space, with blue bars characterizing organs with a z-value of <1.5. The figure shows
 464 z-scores >1.5 for the liver, kidneys, and bladder (red bars). Corresponding organ segmentations are
 465 displayed to the right for the liver (z = 1.9) and the heart (z = 9.9), indicating suboptimal segmentation
 466 results that require manual correction. (B) Error analysis in similarity space for a representative
 467 mesothelioma patient with z-scores >1.5 for the liver, heart, bladder, and lung. Incorrect organ
 468 segmentations are shown to the right for the liver (z = 2.8) and the heart (z = 8.4).



469

470 **Figure 5** (A) Organ segmentation of a hold-out lymphoma test dataset using a training dataset that did
 471 not include splenomegaly cases. (B) Organ segmentation of the same patient following inclusion of 2
 472 (different) lymphoma datasets and model retraining using the expanded training dataset. It can be seen
 473 that the updated model is able to recognise the new image pattern resulting in a correct segmentation of
 474 both the liver and the spleen.
 475

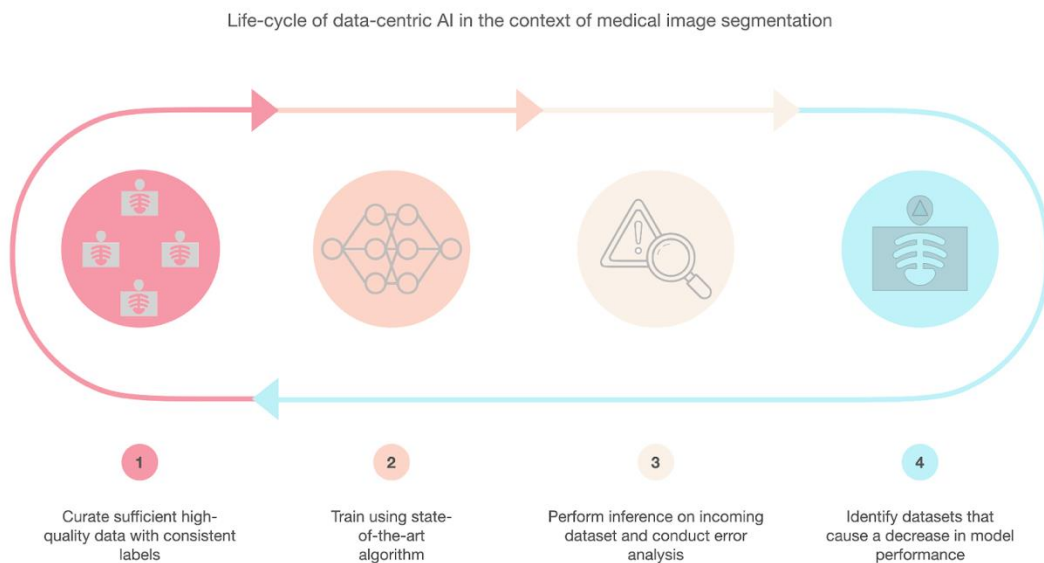


476

477 **Figure 6** (A) Organ segmentation of a hold-out mesothelioma test dataset using a training dataset that did
 478 not include mesothelioma cases. (B) Organ segmentation of the same patient following inclusion of 2
 479 (different) mesothelioma datasets and model retraining using the expanded training dataset. The updated
 480 model recognized the new image pattern resulting in a correct segmentation of the heart.
 481

482 **SUPPLEMENTAL FIGURES**

483

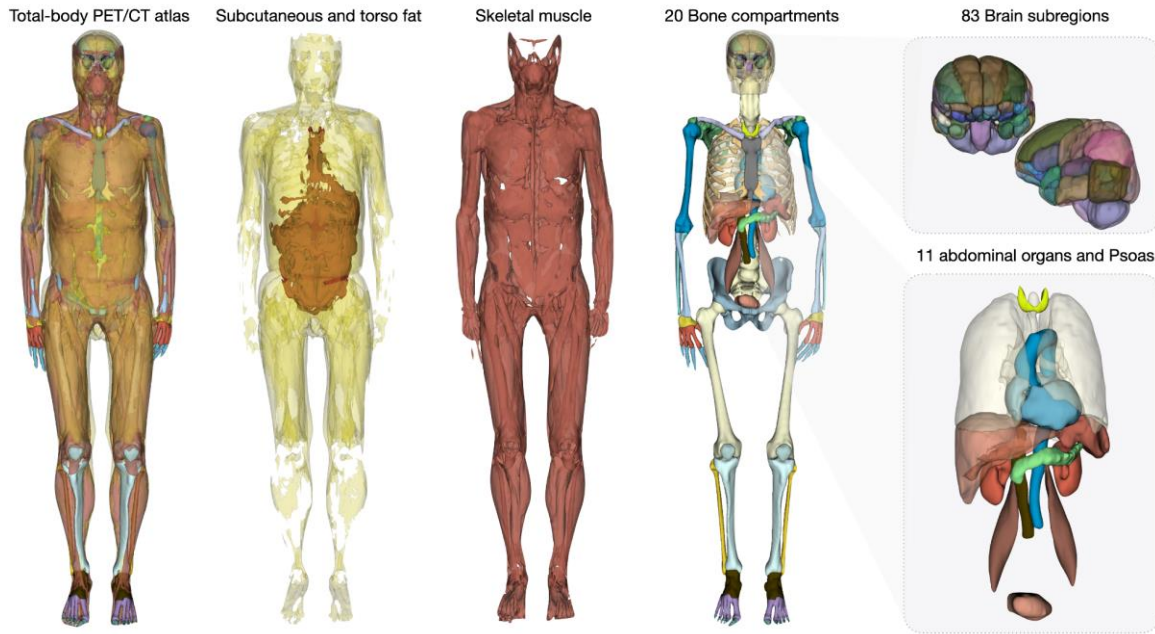


484

485

486 **Supplemental Figure 1.** Life-cycle of data-centric AI approach. The main goal is to identify the datasets
487 that cause a reduction in the model’s performance. Once these datasets are identified, they are added to
488 the original training dataset, and the network is retrained.

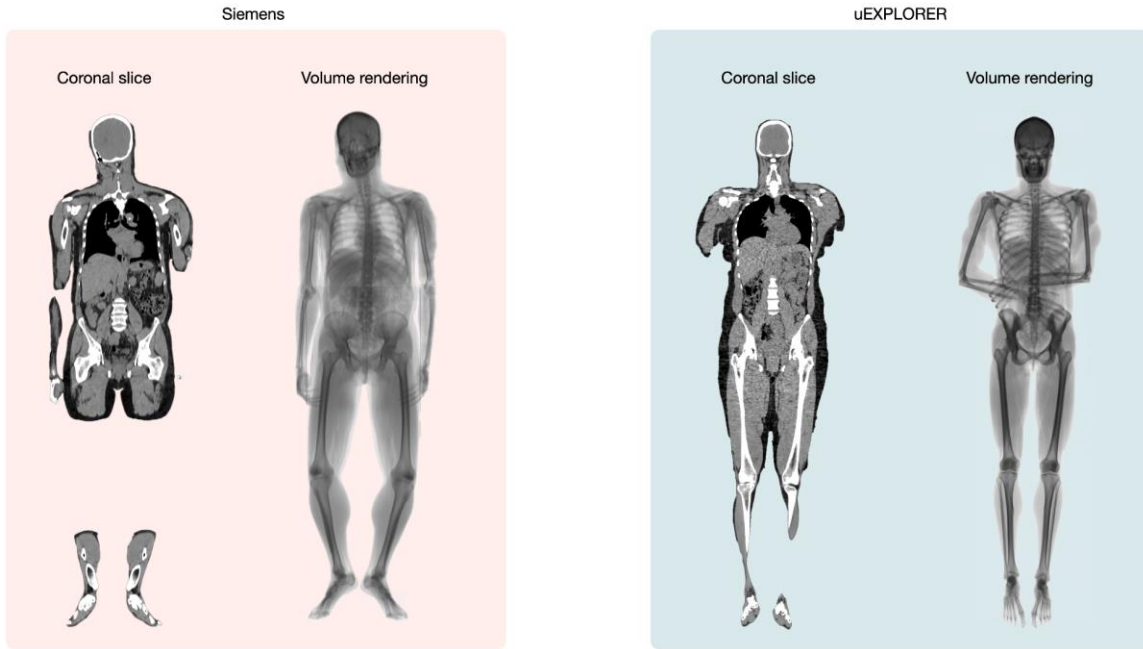
489



490

491 **Supplemental Figure 2.** ^{18}F -FDG total-body PET/CT tissue-map consisting of 120 unique tissues - both
 492 cerebral and non-cerebral structures.

493



494

495 **Supplemental Figure 3.** An example image of a low-dose CT from the Siemens (left) and uEXPLORER (right)

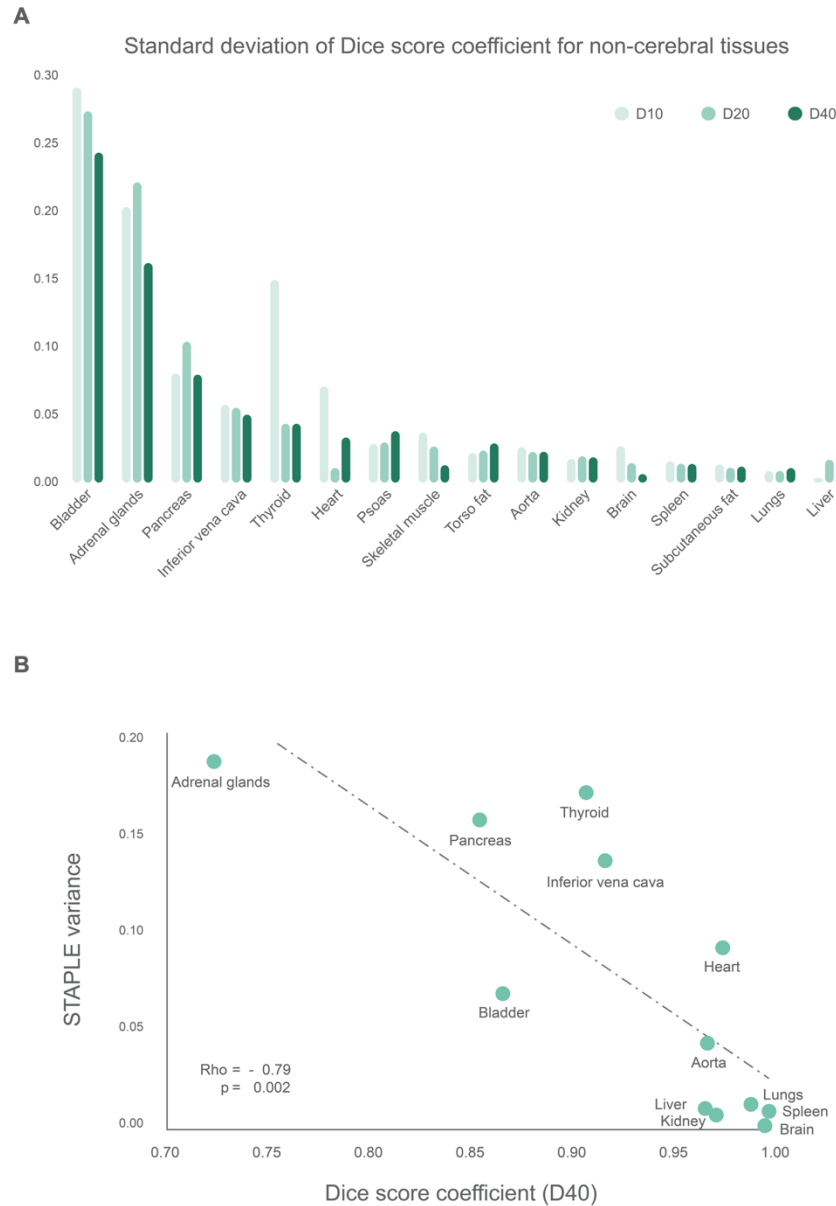
496 system. The two images differ in their noise characteristics and the hand-positions of the subjects.

497 Subjects scanned using the Siemens system position their hands side-by-side (left), while subjects scanned

498 using the uEXPLORER system have their hands crossed across the chest (right).

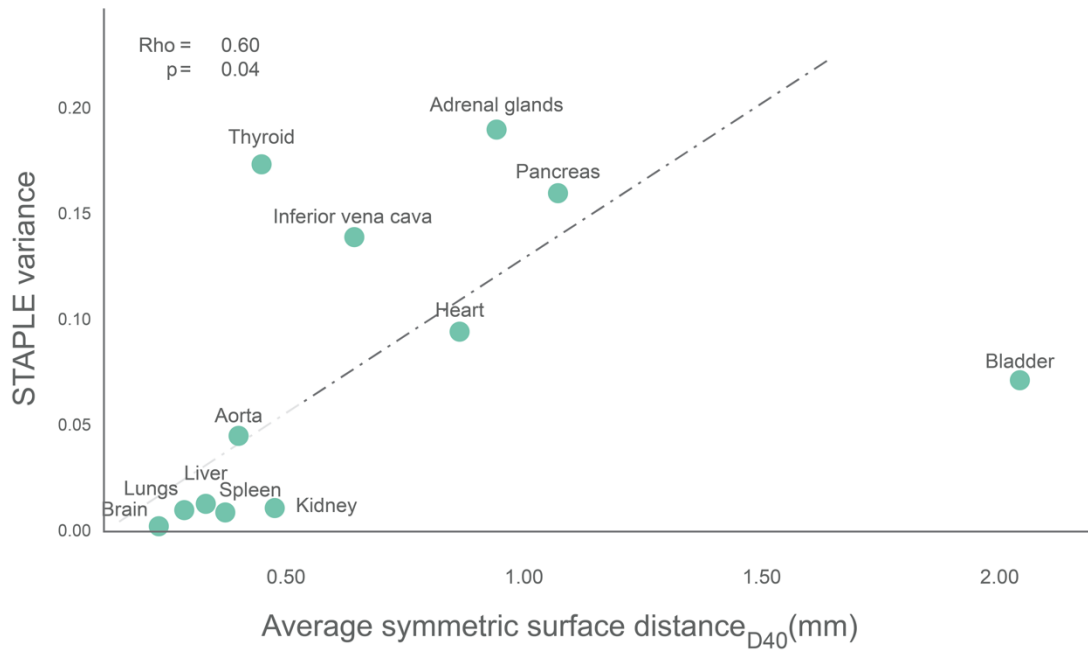
499

500



501

502 **Supplemental Figure 4. (A)** Standard deviation (SD) values of Dice Score Coefficients (DSCs) were
 503 determined for various abdominal organs derived using 10 (D10), 20 (D20) and 40 (D40) training datasets.
 504 The graph indicates that SD values are similar for most organs except for the bladder, adrenal gland,
 505 thyroid and pancreas. **(B)** The correlation plot shows a highly significant negative correlation between the
 506 DSCs obtained from the D40 training dataset and the STAPLE variance. The organs with the lowest DSC
 507 values are also those with the highest SD (adrenal gland, pancreas and thyroid), indicating that
 508 segmentation of these three structures is problematic even for human experts.



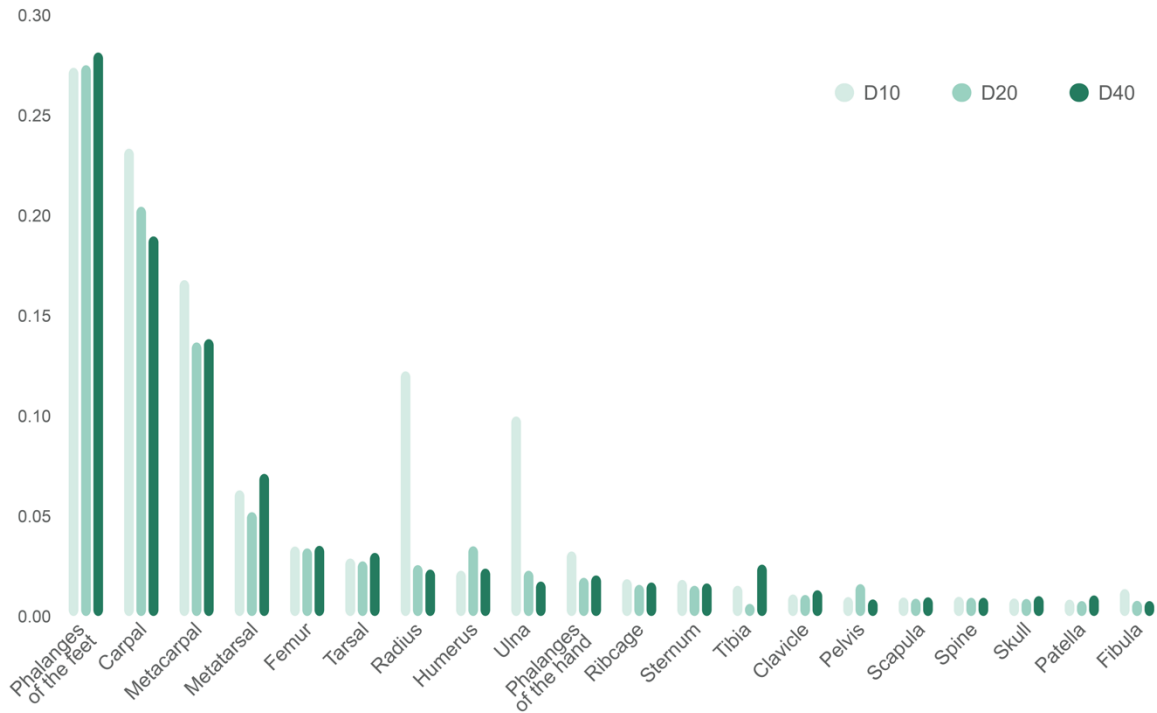
509

510 **Supplemental Figure 5.** The correlation plot shows a significant correlation between the average
 511 symmetric surface distance (ASSD) and the STAPLE variance, indicating that the segmentation of the
 512 bladder, pancreas and the adrenal gland is more challenging than for other organs. It can be seen that the
 513 average misalignment between contours is about twice ($>1\text{mm}$) that of other organs, confirming the
 514 conclusion derived using the DSC metric.

515

516

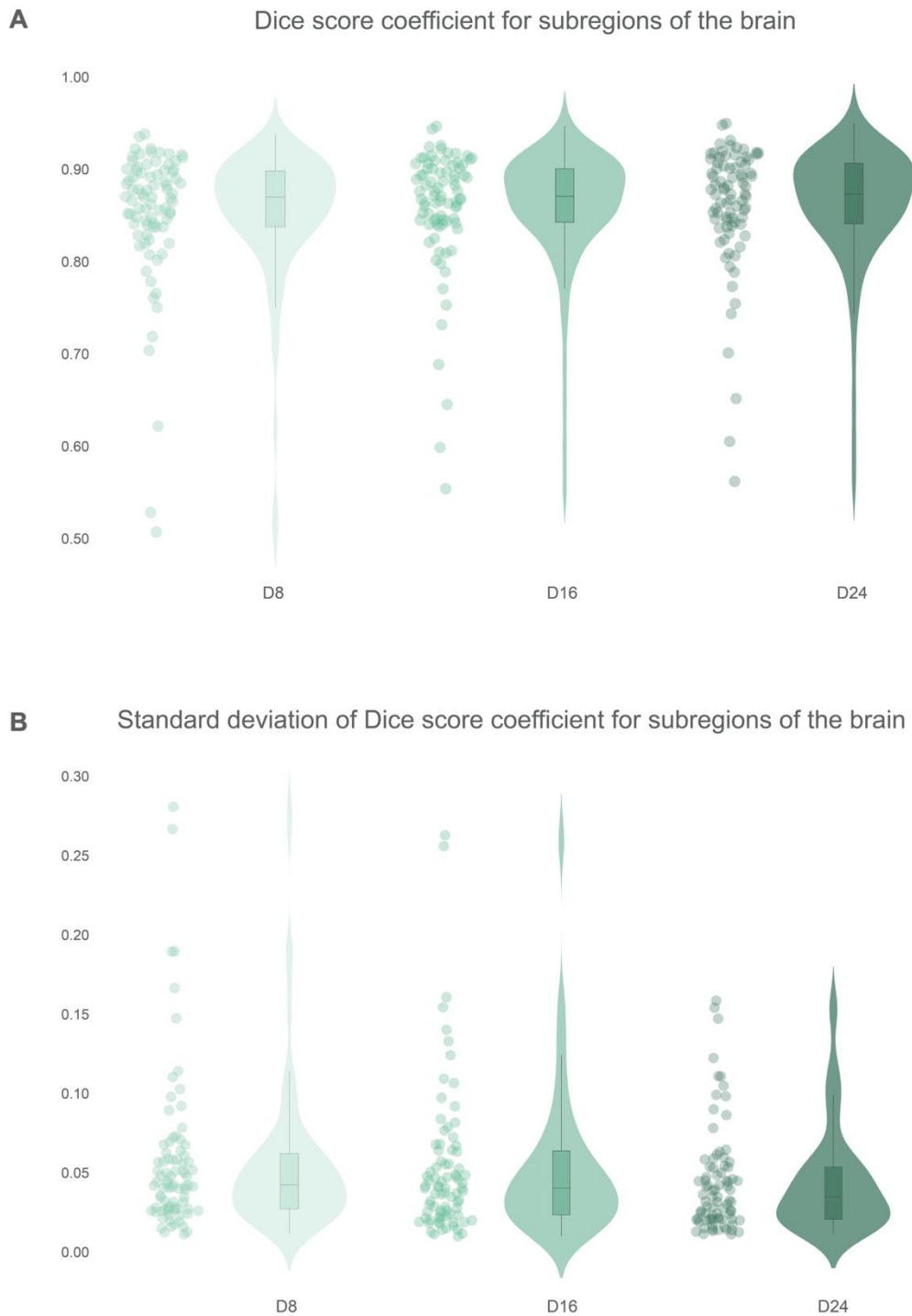
Standard deviation of Dice score coefficient for bones



517

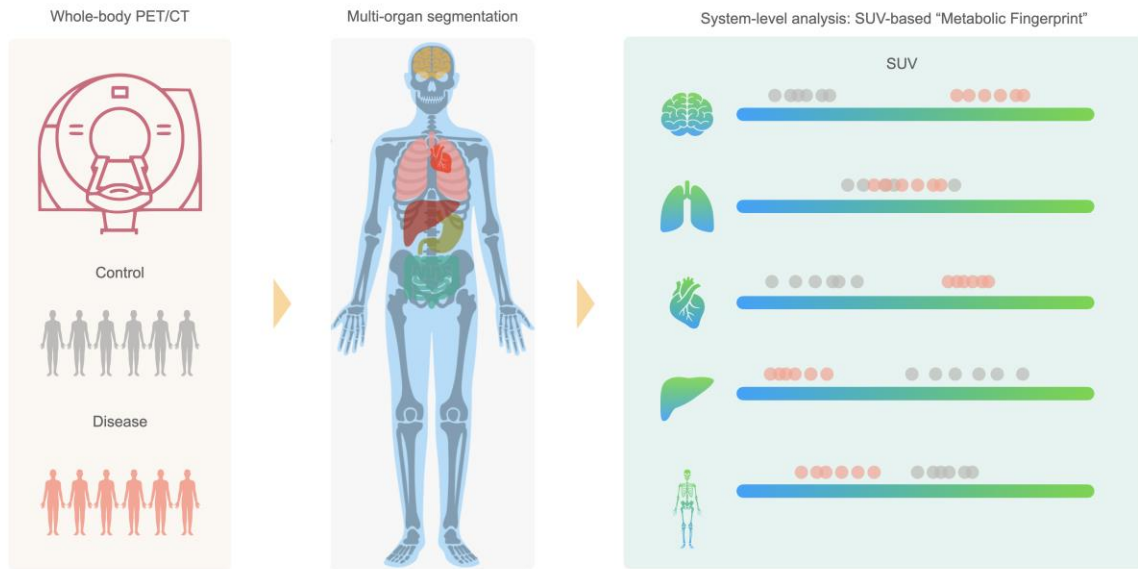
518 **Supplemental Figure 6.** Standard deviation (SD) values of Dice Score Coefficients (DSCs), determined for
519 various bone structures derived using 10 (D10), 20 (D20) and 40 (D40) training data sets. The graph
520 indicates that the SD of the three data sets is similar for most organs, except for the radius and ulna. The
521 variability in segmentation accuracy of these two bone structures is high in case of a small training data.

522



523

524 **Supplemental Figure 7.** Mean (A) and SD (B) values of the DSCs determined for various brain regions
 525 (denoted by dots) derived using 8 (D8), 16 (D16) and 24 (D24) training datasets. One can appreciate that
 526 increase in the training data set size has only a minor effect on segmentation performance, causing a small
 527 decrease in the DSC variance.



528

529 **Supplemental Figure 8.** Concept diagram showing the extraction of a “metabolic fingerprint” from PET/CT
 530 data obtained from a control and disease population. The developed automated multi-organ
 531 segmentation pipeline is used to extract organ SUV values from both groups (control: grey, disease:
 532 orange), allowing the study of systems-level alterations in the pattern of organ metabolic activity as a
 533 consequence of the disease process.

534

535 **SUPPLEMENTAL TABLES**

536 **Supplemental Table 1.** List of cerebral (highlighted in green) and non-cerebral tissues (highlighted in
 537 orange) along with their corresponding label-index in the total-body digital FDG PET/CT tissue-map.

Label-index	Tissues
1	R-Hippocampus
2	L-Hippocampus
3	R-Amygdala
4	L-Amygdala
5	R-Anterior-temporal-lobe-medial-part
6	L-Anterior-temporal-lobe-medial-part
7	R-Anterior-temporal-lobe-lateral-part
8	L-Anterior-temporal-lobe-lateral-part
9	R-Parahippocampal-and-ambient-gyri
10	L-Parahippocampal-and-ambient-gyri
11	R-Superior-temporal-gyrus-posterior-part
12	L-Superior-temporal-gyrus-posterior-part
13	R-Middle-and-inferior-temporal-gyrus
14	L-Middle-and-inferior-temporal-gyrus
15	R-Fusiform-gyrus
16	L-Fusiform-gyrus
17	R-Cerebellum
18	L-Cerebellum
19	Brainstem
20	L-Insula
21	R-Insula
22	L-Lateral-remainder-of-occipital-lobe
23	R-Lateral-remainder-of-occipital-lobe
24	L-Cingulate-gyrus-gyrus-cinguli-anterior-part
25	R-Cingulate-gyrus-gyrus-cinguli-anterior-part
26	L-Cingulate-gyrus-gyrus-cinguli-posterior-part
27	R-Cingulate-gyrus-gyrus-cinguli-posterior-part

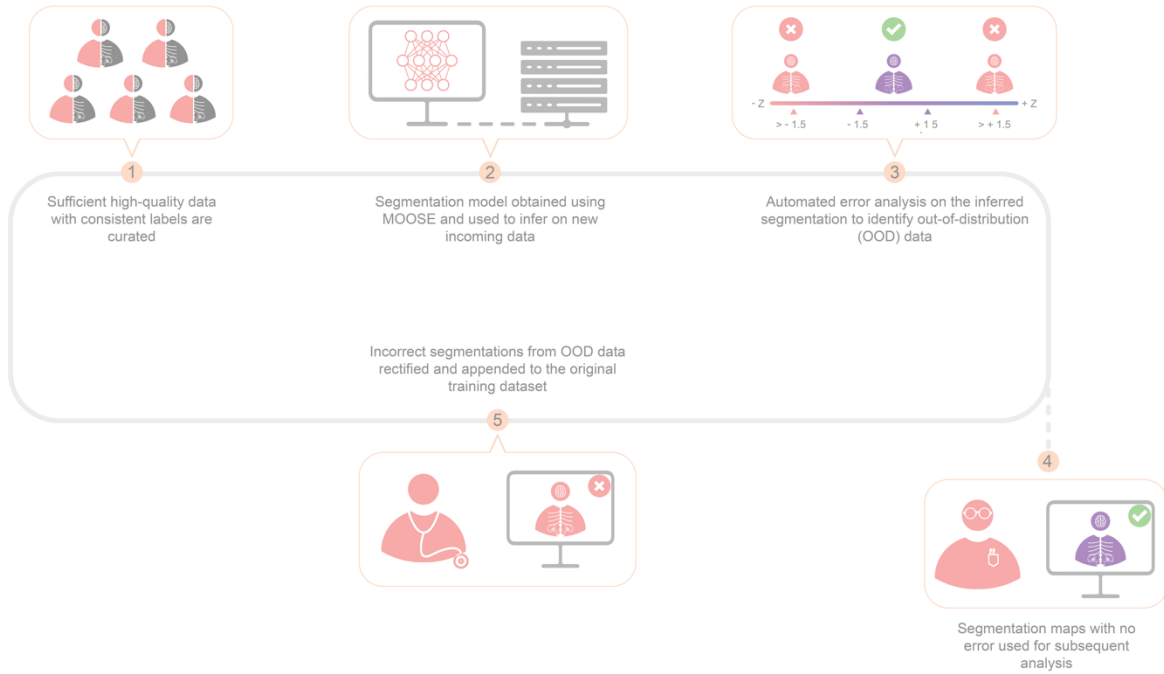
28	L-Middle-frontal-gyrus
29	R-Middle-frontal-gyrus
30	L-Posterior-temporal-lobe
31	R-Posterior-temporal-lobe
32	L-Inferiolateral-remainder-of-parietal-lobe
33	R-Inferiolateral-remainder-of-parietal-lobe
34	L-Caudate-nucleus
35	R-Caudate-nucleus
36	L-Nucleus-accumbens
37	R-Nucleus-accumbens
38	L-Putamen
39	R-Putamen
40	L-Thalamus
41	R-Thalamus
42	L-Pallidum
43	R-Pallidum
44	Corpus-callosum
45	R-Lateral-ventricle-excluding-temporal-horn
46	L-Lateral-ventricle-excluding-temporal-horn
47	R-Lateral-ventricle-temporal-horn
48	L-Lateral-ventricle-temporal-horn
49	Third-ventricle
50	L-Precentral-gyrus
51	R-Precentral-gyrus
52	L-Straight-gyrus
53	R-Straight-gyrus
54	L-Anterior-orbital-gyrus
55	R-Anterior-orbital-gyrus
56	L-Inferior-frontal-gyrus
57	R-Inferior-frontal-gyrus
58	L-Superior-frontal-gyrus
59	R-Superior-frontal-gyrus
60	L-Postcentral-gyrus

61	R-Postcentral-gyrus
62	L-Superior-parietal-gyrus
63	R-Superior-parietal-gyrus
64	L-Lingual-gyrus
65	R-Lingual-gyrus
66	L-Cuneus
67	R-Cuneus
68	L-Medial-orbital-gyrus
69	R-Medial-orbital-gyrus
70	L-Lateral-orbital-gyrus
71	R-Lateral-orbital-gyrus
72	L-Posterior-orbital-gyrus
73	R-Posterior-orbital-gyrus
74	L-Substantia-nigra
75	R-Substantia-nigra
76	L-Subgenua-frontal-cortex
77	R-Subgenua-frontal-cortex
78	L-Subcallosal-area
79	R-Subcallosal-area
80	L-Pre-subgenua-frontal-cortex
81	R-Pre-subgenua-frontal-cortex
82	L-Superior-temporal-gyrus-anterior-part
83	R-Superior-temporal-gyrus-anterior-part
84	Adrenal-glands
85	Aorta
86	Bladder
87	Brain
88	Heart
89	Kidneys
90	Liver
91	Pancreas
92	Spleen

93	Thyroid
94	Inferior vena cava
95	Lung
96	Carpal
97	Clavicle
98	Femur
99	Fibula
100	Humerus
101	Metacarpal
102	Metatarsal
103	Patella
104	Pelvis
105	Phalanges of the hand
106	Radius
107	Ribcage
108	Scapula
109	Skull
110	Spine
111	Sternum
112	Tarsal
113	Tibia
114	Phalanges of the feet
115	Ulna
116	Skeletal-muscle
117	Subcutaneous-fat
118	Torso-fat
119	Psoas
120	Entire Skeleton

539 GRAPHICAL ABSTRACT

540



541

542

543